

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 09-231238

(43)Date of publication of application : 05.09.1997

(51)Int.Cl.

G06F 17/30

(21)Application number : 08-058391

(71)Applicant : OMRON CORP

(22)Date of filing : 20.02.1996

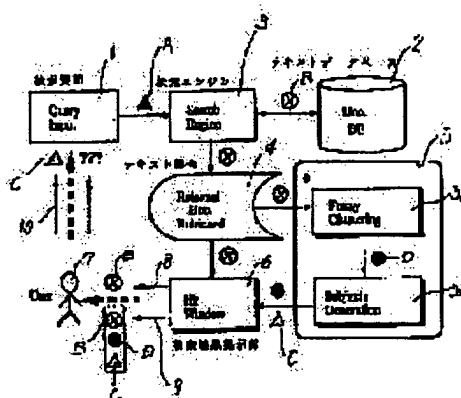
(72)Inventor : GO ATOU
SOGO TAIJI
SAWADA AKIRA

(54) DISPLAY METHOD FOR TEXT RETRIEVAL RESULT AND DEVICE THEREFOR

(57)Abstract:

PROBLEM TO BE SOLVED: To improve both retrieval efficiency and accuracy by dividing a text set into plural groups based on the theme analysis result of every text, generates the theme sort information showing the attribute of every group, and displays these information in every group.

SOLUTION: A retrieval engine 3 expands a retrieval expression based on a prescribed algorithm and extracts a relative text set 4 out of a document data base 2. A fuzzy gathering part 5a of a processing part 5 divides the set 4 into plural groups based on the theme analysis result of every text, and a theme sort information generation part 5b generates the theme sort information showing the attribute of every group. A retrieval result display part 6 processes the acquired information (text set B, centroid D and theme sort information C) in a prescribed display mode and shows them to a user 7. As a result, the document retrieval result can be easily confirmed and the retrieval efficiency and accuracy can be improved owing to prevention of the retrieval omission.



LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

(51)IntCl ⁶	G 0 6 F 17/30	国際記号	庁内整理番号	P I	特許表示箇所
		G 0 6 F	15/403	3 7 0 A	
			15/401	3 1 0 D	
			15/403	3 8 0 E	

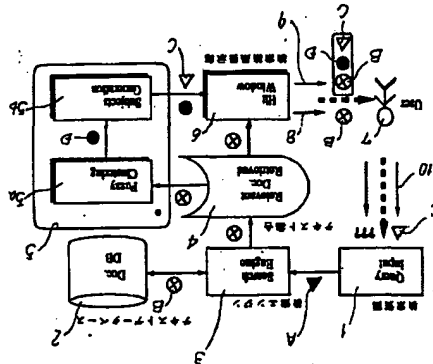
(21)出願番号	特開平9-58381	審査請求	未請求	請求項の数	24	FD (全 19 頁)
(22)出願日	平成8年(1996)2月20日	(71)出願人	000002945	オムロン株式会社		
		(72)発明者	呉 重雄	京都府京都市右京区花園土堂町10番地		
		(72)発明者	山本 隆夫	京都府京都市右京区花園土堂町10番地		
		(72)発明者	十河 大治	ムロン株式会社内		
		(72)発明者	山本 隆夫	京都府京都市右京区花園土堂町10番地		
		(72)発明者	山本 隆夫	ムロン株式会社内		
		(74)代理人	仲野士 飯塚 慎市	京都府京都市右京区花園土堂町10番地		

(54)【発明の名称】 テキスト検索結果表示方法及び装置

(57)【要約】

【課題】 文書検索結果に対する問題を容易として、検索効率の向上、並びに、検索漏れの防止による検索精度の向上を図ることができ、しかも、提示された主題情報にデータ等を如何に効果的に絞り込めるかの指示にもなり、この付加された必要情報を利用して高度な適応検索(Relevance Feedback)を行わせる。

【解決手段】 与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合を各テキストの主題分析結果を用いて自動的に複数個のグループに分割し、該分割により得られた各グループのそれぞれについて、当該グループの属性を表現する主題分群情報を生成し、該生成された各グループの主題分群情報をグループ別に区分して表示する。



【特許請求の範囲】

【請求項1】 与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合を各テキストの主題分析結果を用いて自動的に複数個のグループに分割する分割ステップと、前記分割ステップによって得られた各グループのそれぞれについて、当該グループの属性を表現する主題分群情報を生成する生成ステップと、前記生成ステップで求めた各グループの主題分群情報をグループ別に区分して表示する表示ステップとを具備することを特徴とするテキスト検索結果表示方法。

【請求項2】 与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合を各テキストの主題分析結果を用いて自動的に複数個のグループに分割する分割ステップと、前記分割ステップによって得られた各グループのそれぞれについて、当該グループの属性を表現する主題分群情報を生成する生成ステップと、前記各グループのそれぞれについて、そのグループと前記検索条件との間の適合度を求めるグループ適合度算出ステップと、前記生成ステップで求めた各グループの主題分群情報と、前記適合度算出ステップによって求めた適合度の大きい順に、グループ別に区分して表示する表示ステップとを具備することを特徴とするテキスト検索結果表示方法。

【請求項3】 与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合を各テキストの主題分析結果を用いて自動的に複数個のグループに分割する分割ステップと、前記グループ内の各テキストの内容の分析結果に基づいて、各テキストの当該グループに対する所属度を算出する所属度算出ステップと、前記所属度算出ステップと、前記所属度のグループの中で、テキスト表示対象となるグループを選択するための選択ステップと、前記選択ステップで選択されたグループ内のテキストを前記選択された所属度の順に内容表示する表示ステップとを具備する。

【請求項4】 与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合を各テキストの主題分析結果を用いて自動的に複数個のグループに分割する分割ステップと、前記グループ内の各テキストの内容の分析結果に基づいて、各テキストの当該検索条件に対する所属度を算出する所属度算出ステップと、前記所属度のグループの中で、テキスト表示対象となるグループを選択するための選択ステップと、前記選択ステップで選択されたグループ内のテキストを

前記算出された適合度の順に内容表示する表示ステップとを具備する。

【請求項5】 与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合を各テキストの主題分析結果を用いて自動的に複数個のグループに分割する分割ステップと、前記グループ内の各テキストの内容の分析結果に基づいて、各テキストの当該グループに対する所属度を算出する所属度算出ステップと、前記グループ内の各テキストの内容の分析結果に基づいて、各テキストの当該検索条件に対する所属度を算出する所属度算出ステップと、前記所属度のグループの中で、テキスト表示対象となるグループを選択するための表示対象グループ選択ステップと、

前記各グループ内のテキストを検索条件への適合度順に表示するか、或いは当該グループへの所属度の順に表示するかを選択するための表示順序基準選択手段と、前記表示対象グループ選択ステップで選択されたグループ内のテキストを前記表示順序基準選択手段によって選択された表示順序基準の順に内容表示する表示ステップとを具備する。

【請求項6】 前記分割ステップは、与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合を、ファジィ・クラスティング法を用いて複数個のグループに分割する。

【請求項7】 前記生成ステップにて生成される当該グループの属性を表現する主題分群情報は、当該グループの属性を幾つかのキーワードの組により表すものである。

【請求項8】 前記生成ステップにて生成される当該グループの属性を表現する主題分群情報は、当該グループの属性を幾つかのキーワードの組により表すものである。

【請求項9】 与えられた検索条件に基づいてデータベースを検索することにより得られたテキスト集合の各テキストに対してファジィ・クラスティングを行い、各テキストに各分類カテゴリへの所属度を生成する所属度生成ステップと、

前記生成された所属度を用いて、各文書を1若しくは2以上の分類カテゴリに割り付ける文書割り付けステップと、前記複数個の分類カテゴリの中で、テキスト表示対象

るグループの全体像を把握し易くなり、次の処理のためのグループ選択が非常に容易となる。

【0015】この出願の請求項2（又は請求項14）の発明は、与えられた検索条件に基いてデータベースを検索することにより得られたテキスト集合を各テキストの主成分分析結果を用いて自動的に複数個のグループに分割する分割ステップ（又は手段）と、前記分割ステップ（又は手段）によって得られた各グループのそれぞれについて、当該グループの属性を表現する主成分情報欄を生成する生成ステップ（又は手段）と、前記各グループのそれぞれについて、そのグループと前記検索条件との間の適合度を求めるグループ適合度算出ステップ（又は手段）と、前記生成ステップ（又は手段）で求めた各グループの主成分分析結果を、前記適合度算出ステップによって求めた適合度の大きい順に、グループ別に区分して表示する表示ステップ（又は手段）とを具備する、ことを特徴とするテキスト検索結果表示方法（又は装置）にある。

【0016】そして、この請求項2（又は請求項14）の発明によれば、前記請求項1（又は請求項13）に記載の発明の効果に加えて、検索条件への適合度の順に表示するので、検索目的に合致したグループをグループの内容を確認し易く選択することができる。

【0017】この出願の請求項3（又は請求項15）の発明は、与えられた検索条件に基いてデータベースを検索することにより得られたテキスト集合を各テキストの主成分分析結果を用いて自動的に複数個のグループに分割する分割ステップ（又は手段）と、前記グループ内の各テキストの内容の分析結果に基いて、各テキストの当該グループに対する所屬度を算出する所屬度算出ステップ（又は手段）と、前記複数個のグループの中で、テキスト表示対象となるグループを選択するための選択ステップ（又は手段）と、前記選択ステップ（又は手段）で選択されたグループ内のテキストを前記算出された所屬度の順に内容表示する表示ステップ（又は手段）とを具備する、ことを特徴とするテキスト検索結果表示方法（又は装置）にある。

【0018】そして、この請求項3（又は請求項15）の発明によれば、選択されたグループ内のテキストがグループへの所屬度の順に表示されるので、グループの定数が把握し易くなる。

【0019】この出願の請求項4（又は請求項16）の発明は、与えられた検索条件に基いてデータベースを検索することにより得られたテキスト集合を各テキストの主成分分析結果を用いて自動的に複数個のグループに分割する分割ステップ（又は手段）と、前記グループ内の各テキストの内容の分析結果に基いて、各テキストの当該検索条件に対する適合度を算出する適合度算出ステップ（又は手段）と、前記複数個のグループの中で、テキスト表示対象となるグループを選択するための選択ステップ

（又は手段）と、前記選択ステップ（又は手段）で選択されたグループ内のテキストを前記算出された適合度の順に内容表示する表示ステップ（又は手段）とを具備する、ことを特徴とするテキスト検索結果表示方法（又は装置）にある。

【0020】そして、この請求項4（又は請求項16）の発明によれば、検索条件に通じたグループを選択し、さらにその中のテキストを検索条件の順に表示するので、検索結果をグループ分けしないでテキストを適合度の順に表示する場合よりも、検索条件に対して適切なテキストが順次確実に表示される。

【0021】この出願の請求項5（又は請求項17）の発明は、与えられた検索条件に基いてデータベースを検索することにより得られたテキスト集合を各テキストの主成分分析結果を用いて自動的に複数個のグループに分割する分割ステップ（又は手段）と、前記グループ内の各テキストの内容の分析結果に基いて、各テキストの当該グループに対する所屬度を算出する所屬度算出ステップ（又は手段）と、前記グループ内の各テキストの内容の分析結果に基いて、各テキストの前記検索条件に対する適合度を算出する適合度算出ステップ（又は手段）と、前記複数個のグループの中で、テキスト表示対象となるグループを選択するための選択ステップ（又は手段）と、前記各グループ内の各テキストの適合度順に表示するかを選択するための表示順序基準の所屬度の順に表示するかを選択するための表示順序基準を選択ステップ（又は手段）で選択された表示順序基準の順に表示する表示ステップ（又は手段）とを具備する、ことを特徴とするテキスト検索結果表示方法（又は装置）にある。

【0022】そして、この請求項5（又は請求項17）の発明によれば、ユーザーの目的に応じてテキストの表示順序を変えることができる。

【0023】この出願の請求項6（又は請求項18）に記載の発明は、請求項1（又は請求項13）乃至請求項5（又は請求項17）のいずれかに記載のテキスト検索結果表示方法（又は装置）において、前記分割ステップ（又は手段）は、与えられた検索条件に基いてデータベースを検索することにより得られたテキスト集合を、フuzzy・クラスタリング法を用いて複数個のグループに分割することを特徴とするものである。

【0024】そして、この請求項6（又は請求項18）に記載の発明によれば、ある検索式により探し出された文書集合に対して自動的にoverlapping方式で主内容によるフuzzy分割（主成分分析）が行われるため、検索漏れの防止による検索精度の向上が期待できる。

【0025】この出願の請求項7（又は請求項19）に記載の発明は、請求項1（又は請求項13）又は請求項

2（又は請求項14）に記載のテキスト検索結果表示方法（又は装置）において、前記生成ステップ（又は手段）にて生成される当該グループの属性を表現する主成分情報は、当該グループの属性を備ったのキーワードの組により表すものである、ことを特徴とするものである。

【0026】そして、この請求項7（又は請求項19）に記載の発明によれば、当該グループの属性を備ったのキーワードの組を通して直感的に把握することができ

る。

【0027】この出願の請求項8（又は請求項20）に記載の発明は、請求項1（又は請求項13）又は請求項2（又は請求項14）に記載のテキスト検索結果表示方法（又は装置）において、前記生成ステップ（又は手段）にて生成される当該グループの属性を表現する主成分情報は、当該グループの属性を短い文章により表すものであることを特徴とするものである。

【0028】そして、この請求項8（又は請求項20）に記載の発明によれば、当該グループの属性を短い文章を通して誰にでも判りやすく理解させることができる。

【0029】この出願の請求項9（又は請求項21）に記載の発明は、与えられた検索条件に基いてデータベースを検索することにより得られたテキスト集合の特徴行列に対してフuzzy・クラスタリング法を行い、各グループに各分類カテゴリへの所屬度を生成する所屬度生成ステップ（又は手段）と、前記生成された所屬度を用いて、各文書1若しくは2以上の分類カテゴリグループに割り付けられる文書割り付けステップ（又は手段）と、前記複数個の分類カテゴリグループの中で、テキスト表示対象となる分類カテゴリグループを選択するための分類カテゴリグループ選択ステップ（又は手段）と、前記分類カテゴリグループ選択ステップ（又は手段）で選択された分類カテゴリグループ内のテキストをそのグループに対する適合度の順に内容表示する表示ステップ（又は手段）とを具備する、ことを特徴とするテキスト検索結果表示方法（又は装置）にある。

【0030】そして、この請求項9（又は請求項21）に記載の発明によれば、overlapping手法を用いて各文書1若しくは2以上の分類カテゴリグループに割り付け、その状態にて選択された分類カテゴリグループ内のテキストをそのグループに対する適合度の順に内容表示するため、検索効率の向上、並びに、検索漏れの防止による検索精度の向上を図ることができる。

【0031】この出願の請求項10（又は請求項22）に記載の発明は、前記請求項9（又は請求項21）に記載の発明において、前記文書割り付けステップ（又は手段）は、各文書をその所屬度の上位1個の分類カテゴリグループに割り付ける、ことを特徴とするものである。

【0032】そして、この請求項10（又は請求項22）に記載の発明によれば、請求項9（又は請求項21）に記載の発明の効果に加え、各分類カテゴリグループに

いて常に所属度の高い順に一定数個の文 を表示させることができる。

【0033】この出版の請求項11（又は請求項23）に記載の発明は、前記請求項9（又は請求項21）に記載の発明において、前記文書割り付けステップは、各文書のある属値 α 以上の所属度値を有する分類カテゴリに割り付ける、ことを特徴とするものである。

【0034】そして、この請求項11（又は請求項23）に記載の発明によれば、請求項9（又は請求項21）に記載の発明に加え、各分類カテゴリについて常にある属値 α 以上の所属度値を有する文書を表示させることができる。

【0035】この出版の請求項12（又は請求項24）に記載の発明は、前記請求項9（又は請求項21）に記載の発明において、前記文書割り付けステップは、各文書のカテゴリの確率分布を考慮して分類カテゴリに割り付ける、ことを特徴とするものである。

【0036】そして、この請求項12（又は請求項24）に記載の発明によれば、請求項9（又は請求項21）に記載の発明に加え、各分類カテゴリについてカテゴリーの確率分布を考慮して関連する文書を表示させることができる。

【0037】
【発明の実施の形態】以下に、本発明方法及び装置の好適な実施の形態を添付図面を参照しながら詳細に説明する。

【0038】先ず、本発明方法及び装置が適用されたテキスト検索装置の構成を図1のブロック図により概念的に示す。同図において、1は検索開始時に入力されるべきオリジナル検索質問 (Original Query) や検索絞り込み時に入力されるべきフィードバック検索質問 (Feedback Query) を入力するための検索質問入力部 (Query Inputと記す) であり、具体的には、周知のように、マウスやキーボード等の操作部とそれらの信号を処理する入力用ソフトウェアにより構成される。

【0039】2は検索対象となるテキスト集合に相当するテキスト (文書) データベース (Doc. DBと記す) であり、具体的には、ハードディスクや光ディスク等の大容量記憶媒体に記憶されたテキスト集合やインターネット上に存在するホームページ等のテキスト集合がこれに相当する。

【0040】3はテキスト検索システムの中核に位置する検索エンジン (Search Engineと記す) であり、具体的には、周知のように、前述の検索質問入力部1から入力されるオリジナル検索質問 (Original Query) やフィードバック検索質問 (Feedback Query) を検索条件として所定のアルゴリズムに従って検索式を構築し、前述の文書データベース2から関連するテキスト集合を抽出するソフトウェアがこれに相当する。

【0041】4はこのようにして検索エンジン (Search

Engine) 3により抽出された関連するテキスト集合 (Relevant Doc. Retrievedと記す) であり、後述するように、このテキスト集合4が本発明における加工処理の対象となる。

【0042】5は本発明の要部に相当する加工処理部であり、この加工処理部5はテキスト集合4を各グループの主題分析結果を用いて自動的に検索語のグループに分割する分割手段に相当するフuzzy集合化部 (Fuzzy Clusteringと記す) 5aと、こうして得られた各グループそれぞれについて、当該グループの属性を表現する主題分類情報生成部 (Subject Generationと記す) 5bとを中心として構成されている。

【0043】フuzzy集合化部 (Fuzzy Clustering) 6a及び主題分類情報生成部 (Subject Generation) 5bの作用を図2に概念的に示す。同図において、符号4で示される実線にて囲まれた領域は検索エンジン (Search Engine) 3にて抽出されたテキスト集合 (Retrieved Doc. Retrieved) の全体を表す。

【0044】同様に、符号4a、4b、4cで示される領域にて囲まれた3つの領域はフuzzy集合化部 (Fuzzy Clustering) 6にて分割された3つのグループのそれぞれを表す。

【0045】符号Aで示される黒塗り三角印は、検索開始時に入力されるオリジナル検索質問 (Original Query) を表す。符号Bで示される×入り丸印は、オリジナル検索質問 (Original Query) Aの入力により検索抽出されたテキスト集合4の各構成テキストのそれぞれを表す。

【0046】符号C a、C b、C cで示される3個の白抜き三角印は、グループ4 a、4 b、4 cの属性を表現する主題分類情報 (Group Subject) を表す。尚、これらの主題分類情報C a、C b、C cは検索絞り込みのために用いられ、フィードバック検索質問 (Feedback Query) としても好適なものである。

【0047】符号D a、D b、D cで示される3個の黒塗り丸印は、グループ4 a、4 b、4 cの重心を表す。同様に、符号Dで示される黒塗り四角印は、テキスト集合4の重心を表す。

【0048】図2から明らかなように、フuzzy集合化部 (Fuzzy Clustering) 5 aは、検索の結果得られたテキスト集合4に対して、公知のフuzzyクラスタリング処理を施すことにより、テキスト集合4を複数個 (この例では3個) のグループ4 a、4 b、4 cに分割する。一方、主題分類情報生成部 (Subject Generation) 5 bは、こうして得られた各グループ4 a、4 b、4 cのそれぞれについて、当該グループの属性を表現する主題分類情報C a、C b、C cを生成する。図から明らかなように、このようにして得られる当該グループの属性を表現する主題分類情報C a、C b、C cは、各グループ4 a、4 b、4 cの重心D a、D b、D cととは異なるもの

であり、まさしくそれぞれのグループの属性を端的に表したものである。尚、これらのフuzzy集合化部 (Fuzzy Clustering) 5 a及び主題分類情報生成部 (Subject Generation) 5 bの処理内容については、後に、更に詳しく説明する。

【0049】図1に戻って、6は同様に本発明の要部に相当する検索結果提示部 (Hit Windowと記す) であり、この検索結果提示部6では、前述の経過この検索結果提示部 (Hit Window) 6では、主成分分析結果Cを所定の表示態様に加工したのち、ユーザ (Userと記す) 7に対して提示する、それらの表示態様についても、後に詳細に説明する。

【0050】尚、図1においては、実線により表された本発明による情報の流れと破線により表された従来装置による情報の流れとが同時に示されている。すなわち、従来装置にあつては、検索結果提示部 (Hit Window) 6では、破線矢印8に示されるように、検索の結果得られたテキスト集合Bをそのままユーザ7に提示するのみであり、この場合、テキスト集合Bに示されるテキスト数が多量の場合、目的とするテキストを探し出すのにユーザが不便を来す。これに対して、本発明にあつては、検索結果提示部 (Hit Window) 6では、実線矢印9に示されるように、検索の結果得られたテキスト集合Bのみならず、各分類の重心 (Cluster centroids) D並びに主題分類情報 (Group Subject) Cまでもがユーザ7に提示されることとなるため、特に、この主題分類情報 (Group Subject) Cを手軽に出力して、目的とするテキストを容易に探し出すことが可能となる。すなわち、実線矢印10に示されるように、このようにして得られた主題分類情報C (図2のC1、C2、C3に相当する) をそのままフィードバック検索質問 (Feedback Query) Cとして検索質問入力部 (Query Input) 1に与えれば (図2の実線矢印11に検索質問が分岐 "Query Splitting" する様子を示す)、テキスト集合4を的確に絞り込み、目的とするテキストを容易に探し出すことができ、すなわち高度な適応検索 (Relevance Feedback) を行わせることができるのである。

【0051】次に、以上概念的に説明したテキスト検索装置を、さらにその上面表示態様及びそれを実装するためのデータ処理を中心として、図3以下の図面を参照し詳細に説明する。

【0052】本発明に係るテキスト検索装置におけるデータ処理の全体を図3のゼネラルフローチャートに示す。尚、このゼネラルフローチャートに示される処理は、所定のシステムメニューにおいて、そのメニュー項目のひとつを選択することにより起動される。

【0053】同図において処理が開始されると、検索装置を構成する画像表示部の画面上には所定の表示態様に、より検索画面が表示される (ステップ301)。このようにして表示される検索画面の一例を図4に示す。同図

に示されるように、表示画面は縦長形状のウィンドウW1により構成されており、その上部略3分の1の部分は検索質問入力領域A1とされており、また下部略3分の2の部分は検索結果出力領域A2とされている。

【0054】検索質問入力領域A1内には検索質問入力用のウィンドウW2が設けられており、このウィンドウW2の上側には、入力ガイド文 (Enter Query in plain English) 12が、またその右側には、前述した検索エンジン (Search Engine) 3に対する起動指示を与えるための起動ボタン (図中OKと記す) 13と、検索質問 (Query) を取り消すための取り消しボタン (図中CA、NCELと記す) 14と、システムメニューに対して操作支援等5とが設けられている。

【0055】検索結果出力領域A2内には検索結果出力用のウィンドウW3が設けられており、このウィンドウW3の右側にはスクロールバー16が設けられている。更に、この検索結果出力領域A2の下側には、検索結果としてテキスト全文出力を要求するための全文表示ボタン (図中Full Textと記す) 17と、GBEボタン (図中F1と記す) 18とを設けている。また、検索結果としてテキストの抽出を要求するための抽出ボタン (図中Extractと記す) 21とが設けられている。

【0056】尚、以上の各操作ボタン13、14、15、16、17、18、19、20、21の操作は、カーソルを希望のボタンに移動させた後、マウスのクリック操作等に行われることは言うまでもない。

【0057】そして、入力ガイド文 (Enter Query in plain English) 12に従って、キーボードから検索質問を自然語 (特に、この例では英語) にて、例えば、"I want to know Clinton's political condition." の如くに入力すると、この入力された検索質問22はウィンドウW2内に表示されることとなる。

【0058】この状態において、起動ボタン (図中OKと記す) 13が操作されると、図3に戻って、検索装置が実行され、検索エンジン (Search Engine) 3が起動されて、検索質問に関連するテキスト集合4が文書データベース2より抽出され、この抽出されたテキスト集合の各構成テキストは検索結果22との適合度の高い順にソートされ、そのタイトル23のみがウィンドウW3内に表示される (ステップ302)。また、各テキストのタイトル23の先頭部分には、当該テキストの検索質問に対する適合度を三段階 (「高」、「中」、「低」) に区分して表示する適合度マーク24 a、24 b、24 cが表示される。ここで、黒色塗り済みの丸印にて示される適合度マーク24 aは適合度「高」に、灰色塗り済みの丸印にて示される適合度マーク24 bは適合度

【中】に、白抜きの大印にて示される適合度マーク2 4
oは適合度『低』にそれぞれ対応している。

【0059】以後、図3に示す、システム側において
は文書処理機能の運用を待機する状態となる(ステップ
303)。この状態において、図4の画面に示される分
類化要求ボタン(Gr o u p i n g) 19が操作される
と、本発明の要部である分類化処理が実行される(ステ
ップ306)。

【0060】分類化処理の詳細を図5に示す。同図にお
いて処理が開始されると、所定の案内画面を提示するこ
とにより、分類グループ数gの指定を待機する状態とな
る(ステップ501)。この状態において、分類グルー
プ数gの指定(この例では『5』)が完了すると、本発
明の特微部分である文 特徴量の抽出処理(ステップ5
02)、フuzzy・クラスタリング処理(Fuzzy Clust
eringと記す)(ステップ503)、及び主題分類情報の
生成処理(ステップ504)が順に実行される。

【0061】文書特徴量の抽出処理(ステップ502)
では、次のようにして、文書抽象化と文書特徴ベクトル
の生成が行われる。文書は重み付けられた語の集合(語
を構成要素とするベクトル)によって表され、文書の集
合は語を構成要素とする行列として表される。そのた
め、各文書の特徴となる単語(重要語)を自動的に切り
出し、単語の出現を次元mとし、各要素が文書単位の単
語の出現頻度(出現回数)に比例するようベクトル表現F iを用い
ることによって、文書は数1の如くに抽象化される。

【0062】

【数1】

$$F_i = (f_{i1}, f_{i2}, \dots, f_{im})$$
$$f_{ij} = \text{文書 } i \text{ の特徴ベクトル } f_i \text{ の } j \text{ 成分 (} f_{ij} \text{ : 文書 } i \text{ の } j \text{ 成分に対する重み (原素、或は他の評価値))}$$

数1

文 ベクトル集合の例を表1に示す。この例では、文書
集合の構成文書(F 1, F 2, F 3...)のそれぞれに含
まれる重要語(C l i n t o n, S i n g a p o r e, C h i n a...)の重み
(例えば、値域)が示されている。

【0063】

【表1】

	C l i n t o n	S i n g a p o r e	C h i n a	...
F1	0.8	0.4	0.0	
F2	0.6	0.7	0.0	
F3	0.5	0.0	0.7	
...				

文書ベクトル集合の例

数1

表1に示される文書ベクトル集合を文書空間に展開した
例を図6に示す。この例では、前述の重要語(C l i n t o n,
S i n g a p o r e, C h i n a...)を座標軸とする文書空間に文書集
合の各構成文書(F 1, F 2, F 3...)が展開されてい
る。

【0064】続くフuzzy・クラスタリング処理(ステ
ップ503)では、検索結果としての文書集合の特徴行
列に対し、公知のFCM法を用いてフuzzy・クラスタ
リングを行うことにより、次の2種類の分類情報(V
c, U i)が生成される。

【0065】1) 各分類の代表文書特徴ベクトルV c
【数2】

$$V_c = (vc1, vc2, \dots, vc_m)$$
$$V_g : \text{分類グループ } g \text{ の代表文書ベクトル}$$
$$V_c j : \text{分類グループ } g \text{ の } j \text{ 成分に対する重み}$$

数2

2) 各文書の各分類カテゴリへの所属度U i

【数3】

$$U_i = (u_{i1}, u_{i2}, \dots, u_{ig})$$
$$U_i : \text{文書 } i \text{ の各分類グループへの所属度ベクトル}$$
$$u_{ij} : \text{文書 } i \text{ の } j \text{ 分類グループへの所属度}$$

数3

文書分類所属度の例を表2に示す。この例では、各文書
の所属度(U 1, U 2, U 3...)が各分類グループ(G
1, G 2, G 3...)毎に示されている。

【0066】

【表2】

	G1	G2	...	Gg
U1	0.7	0.2	...	
U2	0.8	0.1	...	
U3	0.3	0.6	...	
...				

文書分類所属度の例

数2

続く分類主題情報の生成処理(ステップ504)では、
次の2種類の方式により、分類主題情報の生成が行われ
る。

【0067】1) キーワード方式
このキーワード方式は、各分類グループの主題を導つか
このキーワードの組み合わせにより表現する方式であり、
その際、キーワードの抽出には次の2種類の方式が考
えられる。第1の方式は、該当分類の代表文書ベクトル
V cにおける単語の重みの高い単語を順番にk個抽出し
てそれらの単語をそのグループの主題を表す情報として
用いるものである。第2の方式は、該当分類の文書集合
に対して所属度の高い順にr個の文書ベクトルを抽出
し、そのr個の文書ベクトル集合において出現文書数の
高いものから順にk個の単語を抽出して、そのグループ
の主題情報を表す情報として用いるものである。

【0068】2) テキスト方式

このテキスト方式では、上記のキーワード方式で主題情
報を生成するために抽出されたr個の文書の先頭段落の
テキスト(タイトルを含む)に対し、キーワード方式で
得られたキーワード主題情報を利用して文単位で文字列
照合によりそれらのキーワードを最も多く所有するテキ
ストを抽出し、そのテキスト文をそのグループの主題情
報として用いるものである。

【0069】このようにして得られた各グループの主題
情報、すなわち分類主題情報(前述のキーワード群又は
タイトル等)は、後述するように、所定の提示順番に
ユーザに提示されることとなる。ここで、検索された
文書1の検索結果間に対する適合度をR i、分類グルー
プの検索結果間の適合度をG R cとすると、両者間には数4
の関係が成立する。

【0070】

【数4】

$$G R c = \sum_{i=1}^r R i / r c$$
$$r c : \text{グループ } c \text{ に対して所属度の高い順に抽出された文書数}$$
$$R i : \text{グループ } c \text{ に対して抽出された } r c \text{ 個の文書集合内の } i \text{ 文書の適合度}$$

数4

ここで、数4に示された、グループcに対して所属度の
高い順に抽出された文書数r c (c=1, ..., g)は、
分類数の求め方を図7のフローチャートに示す。同図
において、処理が開始されると、r cの初期化(r c=
0)を行ったのち(ステップ701)、文書1の所属度
の行データU i1に対して最大の所属度が求められ(ステ
ップ702)、その最大値と対応しているグループcの
メンバ数r cが加算され(ステップ703)、以上の処
理(ステップ702, 703)がiを+1づつ加算しつ
つ(ステップ704)、その加算値がi=n(文書数)
となるまで(ステップ705YES)繰り返されて、そ
の結果r cの値が最終的に求められることとなる。
【0071】このようにして、分類主題情報の生成(提
示順番の決定を含む)が完了すると(ステップ50
4)、求められた分類主題情報を用いた検索結果の動的
表示処理が開始される(ステップ505)。

【0072】検索結果の動的表示処理の詳細を図8のフ
ローチャートに示す。同図において処理が開始される
と、検索装置を構成する画像表示器の画面上に設定され
た検索結果出力領域A 2は、図9又は図10に示される
ように、上下に2分割され、これにより分類主題情報
を用いたウィンドウ (Subject Window) W4と検索結果出力
ウィンドウ (Hit Window) W5とが得られる。そして、
分類主題情報表示ウィンドウ W5 (Subject Window) W4
において、所定の表示態様により、各分類主題情報の提
示が行われる(ステップ801)。前述したように、こ
の各分類主題情報の提示は、キーワード方式とテキス
ト方式とで行われる。

【0073】キーワード方式による表示画面の一例を図
9に示す。尚、この例では、検索されたテキスト集合が
5個の分類グループに分割されている。同図に示される

ように、主題分類情報表示用ウィンドウ (Subject Window) W4 内には、その左端部に沿うようにして、分類グループ番号「1」～分級グループ番号「5」に対応する5個のグループボタン25～29が上下一列に配置され、お、それらのグループボタン25～29の右側には、当該分級グループの主題を的確に表すキーワード群30～34が配列されている。この例では、分級グループ番号「1」に対応するグループボタン25の右側には、キーワード群30として、"SINGAPORE:CHINESE:POLITICS"が配列されており、分級グループ番号「2」に対応するグループボタン26の右側には、キーワード群31として、"DALAILAMA:MEET:CHINA:TIBET"が配列されており、分級グループ番号「3」に対応するグループボタン27の右側には、キーワード群32として、"WORLD:LEADER:GOVERNMENT:OFFICIAL"が配列されており、分級グループ番号「4」に対応するグループボタン28の右側には、キーワード群33として、"WIKI:WIKI:WIKI:WIKI"が配列されており、分級グループ番号「5」に対応するグループボタン29の右側には、キーワード群34として、"QUESTION:CHARACTER:PEOPLE:POLITICS"が配列されている。

【0074】また、これらの主題分類情報に、先に求められた提示順番に従い、検索質問 (Query) との適合度の高いものから順に配列されている。すなわち、この例では、分級グループ「1」にて検索される主題が最も検索質問との適合度が高く、分級グループ「5」にて検索される主題が最も検索質問との適合度が低いこととなる。従って、ユーザーは主題分類情報表示用ウィンドウ (Subject Window) W4 内の表示順番から、自分の探している情報に最も近い分級グループを容易に知ることで、しかもそれぞれの内容の正確性に裏づけされる主題が最も検索質問との適合度が高いこととなる。従って、ユーザーは主題分類情報表示用ウィンドウ (Subject Window) W4 内の表示順番から、自分の探している情報に最も近い分級グループを容易に知ることで、しかもそれぞれの内容の正確性に裏づけされる主題が最も検索質問との適合度が高いこととなる。従って、ユーザーは主題分類情報表示用ウィンドウ (Subject Window) W4 内の表示順番から、自分の探している情報に最も近い分級グループを容易に知ることで、しかもそれぞれの内容の正確性に裏づけされる主題が最も検索質問との適合度が高いこととなる。

【0075】テキスト方式による表示画面の一例を図10に示す。尚、この例でも、検索されたテキスト集合が5個の分級グループに分割されている。図10に示されるように、主題分類情報表示用ウィンドウ (Subject Window) W4 内には、その左端部に沿うようにして、分級グループ番号「1」～分級グループ番号「5」に対応する5個のグループボタン25～29が上下一列に配置されており、それらのグループボタン25～29の右側には、当該分級グループの主題を的確に表す短いテキスト文35～39が配列されている。この例では、分級グループ番号「1」に対応するグループボタン25の右側には、キーワード群35として、"Clinton Protest Singapore

ore Caning, Mulls Response"が表示されており、分級グループ番号「2」に対応するグループボタン26の右側には、キーワード群36として、"Clinton Meets With Dalai Lama"が表示されており、分級グループ番号「3」に対応するグループボタン27の右側には、キーワード群37として、"Indian Leader Meets Clinton"が表示されており、分級グループ番号「4」に対応するグループボタン28の右側には、キーワード群38として、"Nixon Had Living Will"が表示されており、分級グループ番号「5」に対応するグループボタン29の右側には、キーワード群39として、"Clinton News Conference Text"が表示されている。

【0076】また、これらの主題分類情報についても、先に求められた提示順番に従い、検索質問 (Query) との適合度の高いものから順に配列されている。すなわち、この例では、分級グループ「1」にて検索される主題が最も検索質問との適合度が高く、分級グループ「5」にて検索される主題が最も検索質問との適合度が低いこととなる。従って、ユーザーは主題分類情報表示用ウィンドウ (Subject Window) W4 内の表示順番から、自分の探している情報に最も近い分級グループを容易に知ることで、しかもそれぞれの内容の正確性に裏づけされる主題が最も検索質問との適合度が高いこととなる。従って、ユーザーは主題分類情報表示用ウィンドウ (Subject Window) W4 内の表示順番から、自分の探している情報に最も近い分級グループを容易に知ることで、しかもそれぞれの内容の正確性に裏づけされる主題が最も検索質問との適合度が高いこととなる。

【0077】次に、先に説明したフレイ・クラスタリングにより得られた各文書の各分級グループへの所属度U1を用いた、検索結果の最終表示のための処理について詳細に説明する。尚、この例では、分級結果の最終表示のために3種類の処理が用意されており、これらの処理は図9又は図10に示される画面において、グループボタン25～29のいずれか一つを選択することにより起動される (ステップ802)。

【0078】先に説明したように、本発明では検索結果としての文書集合の検索結果に対し、FCM法を用いて、各分級グループへの所属度U1が求められている。今仮に、5個の文書 (001, 002, 003, 004, 005) が存在し、それらの文書のそれぞれについて3個の分級グループ (カテゴリ1, カテゴリ2, カテゴリ3) のそれぞれに対する所属度が表3の通りであると想定する。

【表3】
【0079】

文書番号	カテゴリ1	カテゴリ2	カテゴリ3
001	0.50	0.30	0.20
002	0.60	0.10	0.30
003	0.10	0.80	0.10
004	0.25	0.34	0.41
005	0.35	0.10	0.55

振り付けの説明のための数値例

表3

ログラムの一例を図11に示す。図11において処理が開始されると、k値の設定処理 (ステップ1101) 及び1, c, Ncの初期化処理 (ステップ1102) を実行した後、文書iの所属度データiに対するソート処理 (ステップ1103)、最大所属度データ値から順にk個のグループ番号を抽出する処理 (ステップ1104)、及び該当するk個のグループに文書iを登録すると同時にメンバー値を加算する処理 (ステップ1105) が、文書番号iがnになるまで繰り返され (ステップ1106)、文書番号iがnに達すると各グループ毎の文書割り付け結果を出力して処理が終了 (ステップ1107) する。

【0082】(2) ある閾値α以上の所属度値を有する分級カテゴリに割り付ける場合
この表示処理にあつては、各文書 (001～005) はある閾値α以上の所属度値を有する分級カテゴリに割り付けられる。ここで、αとしては、例えば1/g (g: 分級数) とすることが考えられる。表3に示される例では、g=3, α=0.33となるため、文書 (001) については所属度値が0.33以上であるカテゴリ1に、文書 (002) については同様なる理由でカテゴリ1に、文書 (003) については同様なる理由でカテゴリ2に、文書 (004) については同様なる理由でカテゴリ2とカテゴリ3に、文書 (005) については同様なる理由でカテゴリ1とカテゴリ3に割り付けられる。これを分級カテゴリ1とカテゴリ3に割り付けると、
カテゴリ1G1= (001, 002, 005) ; N1=3
カテゴリ2G2= (003, 004) ; N2=2
カテゴリ3G3= (004, 005) ; N3=2

となり、分級グループG1に含まれる文書数N1は3個、分級グループG2に含まれる文書数N2は2個、分級グループG3に含まれる文書数N3は2個とされる。そして、このようにして各カテゴリに属することとされた文書が、後に詳細に説明するように、グループ番号の指定と共に検索結果出力用ウィンドウ (HitWindow) W5内に表示されることとなる。

【0083】以上の表示処理 (2) を実行するためのプログラムの一例を図12に示す。図12において処理が開始されると、α値の設定処理 (ステップ1201) 及び1, c, Ncの初期化処理 (ステップ1202) を実行した後、文書iの所属度データiに対するu1>α

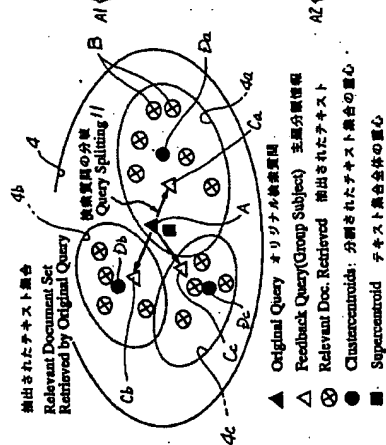
フローチャートである。

【符号の説明】

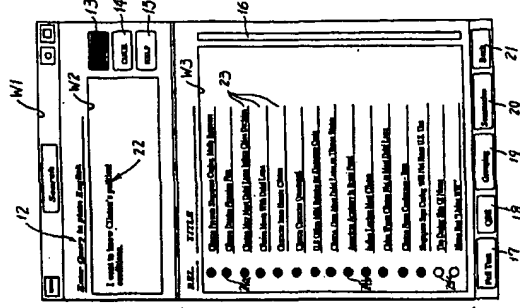
- 1 検索質問入力部
- 2 文データベース
- 3 検索エンジン
- 4 抽出された関連テキスト集合
- 4a, 4b, 4c 分類グループ
- 5 加工処理部
- 5a フラジ集合理化部
- 5b 主題分類情報生成部
- 6 検索結果提示部
- 7 ユーザー
- 12 入力ガイド文
- 13 起動ボタン
- 14 取り換えボタン
- 15 ヘルプボタン
- 16 スクロールバー
- 17 全文要求ボタン
- 18 QBEボタン
- 19 分類要求ボタン
- 20 抄録要求ボタン
- 21 復帰ボタン
- 22 検索質問
- 23 テキスト集合を構成する各テキストのタイトル

- 24a, 24b, 24c 適合度マーク
- 25~29 グループボタン
- 30~34 キーワード群
- 35~39 テキスト文
- 40~44 テキスト文
- 45, 46 スクロールバー
- 49 分類グループ数の表示
- 47a, 47b, 47c グループ毎の適合度マーク
- 48a, 48b, 48c 検索質問に対する適合度マーク
- 49 分類グループ数の表示
- 50 検索質問適合度順指定ボタン
- 51 グループ適合度順指定ボタン
- A オリジナル検索質問
- B 抽出された各構成テキスト
- Ca, Cb, Cc 主題分類情報
- Da, Db, Dc グループの重心
- A1 検索質問入力領域
- A2 検索結果出力領域
- W1 検索質問入力力のウィンドウ
- W2 検索結果出力力のウィンドウ
- W3 主題分類情報表示用ウィンドウ
- W4 主題分類情報表示用ウィンドウ
- W5 検索結果出力用ウィンドウ

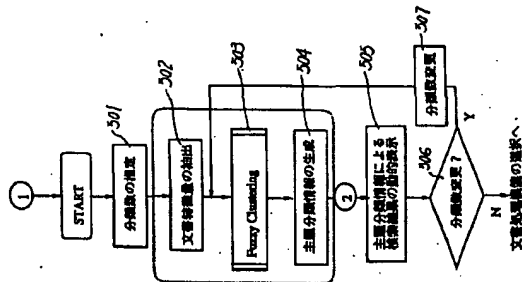
【図2】



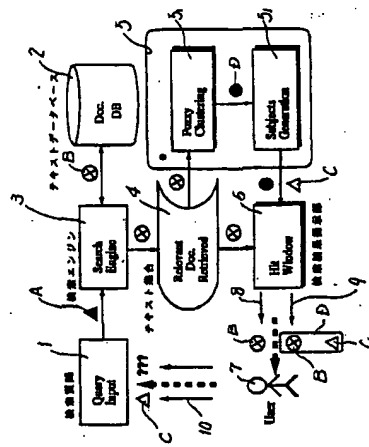
【図4】



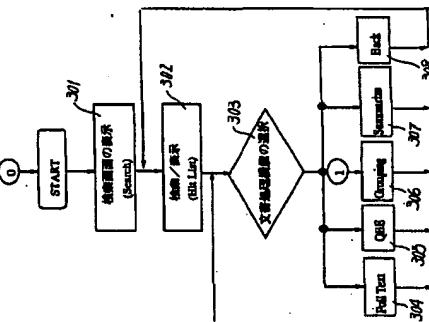
【図5】



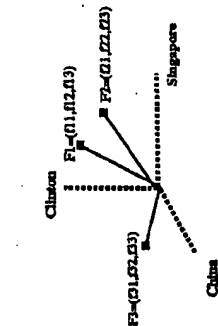
【図1】



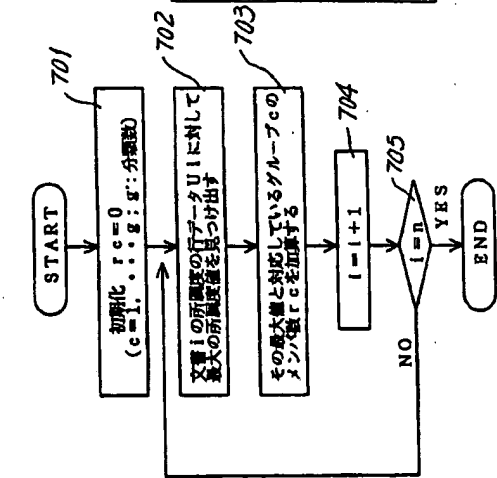
【図3】



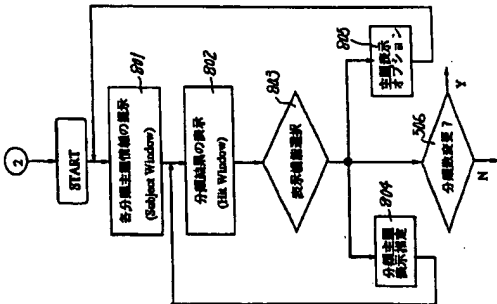
【図6】



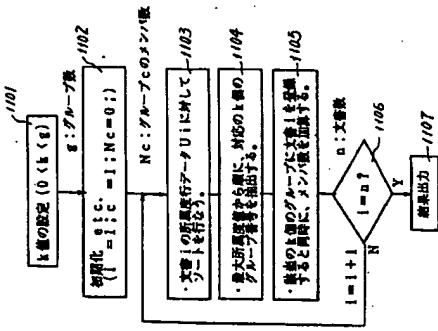
【図 7】



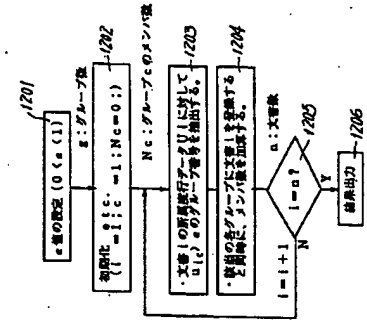
【図 8】



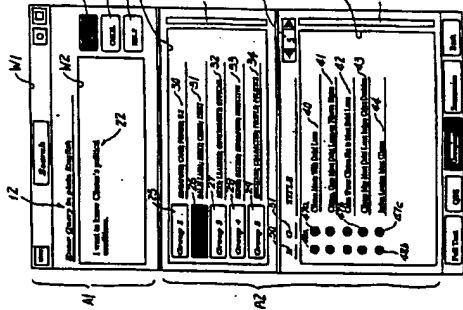
【図 11】



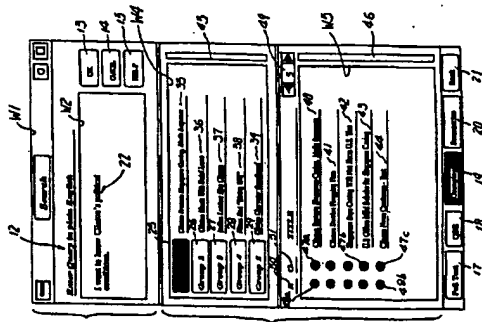
【図 12】



【図 9】



【図 10】



【図13】

